TREE-OF-TABLE: UNLEASHING THE POWER OF LLMS FOR ENHANCED LARGE-SCALE TABLE UNDERSTAND ING

Anonymous authors

Paper under double-blind review

ABSTRACT

The ubiquity and value of tables as semi-structured data across various domains necessitate advanced methods for understanding their complexity and vast amounts of information. Despite the impressive capabilities of large language models (LLMs) in advancing the natural language understanding frontier, their application to large-scale tabular data presents significant challenges, specifically regarding table size and complex intricate relationships. Existing works have shown promise with small-scale tables but often flounder when tasked with the complex reasoning required by larger, interconnected tables found in real-world scenarios. To address this gap, we introduce "Tree-of-Table", a novel approach designed to enhance LLMs' reasoning capabilities over large and complex tables. Our method employs Table Condensation and Decomposition to distill and reorganize relevant data into a manageable format, followed by the construction of a hierarchical Table-Tree that facilitates tree-structured reasoning. Through a meticulous Table-Tree Execution process, we systematically unravel the tree-structured reasoning chain to derive the solutions. Experiments across diverse datasets, including WikiTQ, Table-Fact, FeTaQA, and BIRD, demonstrate that Tree-of-Table sets a new benchmark with superior performance, showcasing remarkable efficiency and generalization capabilities in large-scale table reasoning.

032

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

033 Tables, as a pivotal form of semi-structured data, ubiquitously underpin numerous aspects of daily life 034 and professional domains, ranging from open data repositories and web pages to critical applications in financial analysis, risk management, health monitoring, and business reporting (Cafarella et al., 2008). The advent of large language models (LLMs) (OpenAI, 2023; Yao et al., 2023b; Chen, 2023; Jiang et al., 2023; Imani et al., 2023; Anil et al., 2023; Valmeekam et al., 2022) has opened new 037 vistas for understanding and reasoning with tabular data, marking a significant stride in the realm of natural language understanding (Nahid & Rafiei, 2024; Chen et al., 2024; Sui et al., 2024b;a; Ye et al., 2023; Cheng et al., 2022; Jin & Lu, 2023). This intersection is not only instrumental in 040 enhancing the comprehension of tables but also vital for powering a plethora of downstream tasks 041 such as table-based fact verification (Chen et al., 2019) and question answering (Li et al., 2024). 042 Unlike their unstructured text counterparts, tables provide a dense, structured format through the 043 interaction of rows and columns, offering a rich source of information. However, the same structural 044 characteristics pose unique challenges for language models, as they necessitate advanced levels of reasoning over both the textual and numerical data contained within. Given the increasing reliance on tables for data representation and the complexities involved in their interpretation, investigating the 046 integration of LLMs for improved large-scale table understanding has emerged as an essential and 047 compelling research avenue, drawing heightened interest from the global academic and industrial 048 research communities. 049

Existing methods for table understanding have shown substantial progress in comprehending small scale tables (Cheng et al., 2022; Ye et al., 2023; Wang et al., 2024). However, these approaches
 often falter when applied to the more complex and larger tables frequently encountered in real-world
 scenarios. This gap between academia and practical applications stems from a variety of limitations
 inherent in current methodologies. One significant challenge is the limited contextual capacity of



Figure 1: Comparison of (a) Generic Reasoning, (b) Chain-of-Table (Wang et al., 2024), and the proposed (c) Tree-of-Table methods when confronted with large-scale relational tables. Generic Reasoning often struggles with the increased context and complexity, leading to inefficient processing and potential loss of critical information. Chain-of-Table, while more structured with linear thought chain, still faces challenges with the scale and intricacy of data. In contrast, Tree-of-Table showcases a structured and hierarchical reasoning process that adeptly handles large-scale tables, significantly enhancing comprehension and efficiency compared to previous methods, particularly in managing the complexity of expansive tabular data.

081 082

today's language models. As tables increase in size, the amount of information that must be processed
 and understood grows exponentially due to the intricate interactions between rows and columns. This
 complexity makes it difficult for models to capture and reason about all the necessary information in
 one go, significantly impeding their understanding capabilities. When faced with complex question answering logic that spans lengthy chains, pinpointing, extracting, and comprehending key table
 information becomes an immense challenge.

To address these issues, two main approaches have generally been adopted, as shown in Figure 1. The first involves using only the schema information of tables and employing program-aided methods, 090 such as generating SQL-based answers from questions (Rajkumar et al., 2022b;a; Shi et al., 2020; 091 Pönighaus, 1995; Katsogiannis-Meimarakis & Koutrika, 2023). While this approach avoids directly 092 inputting entire tables, the resulting SQL statements can be lengthy and prone to errors, leading to 093 suboptimal performance. The second strategy involves decomposing tables into multiple sub-tables. 094 Methods like Dater (Ye et al., 2023) attempt to manage larger tables by initially inputting the entire table before breaking it down, which is impractical. The Chain-of-Table (Wang et al., 2024) draws 096 inspiration from the chain-of-thought principle (Wei et al., 2022), performing implicit sub-table extraction. Yet, even this approach is limited to understanding smaller tables. Additionally, traditional 098 table understanding datasets like WikiTQ (Pasupat & Liang, 2015) and TableFact (Chen et al., 2019), which are relatively small, severely restrict the exploration of large-scale table understanding. 099 Fortunately, the introduction of the BIRD (Li et al., 2024) dataset, considered the largest and most 100 complex table understanding dataset to date, highlights the pressing need for improvements. Despite 101 this, due to the reasons mentioned above, existing large language models still exhibit low accuracy on 102 comprehensive, large-scale table datasets like BIRD (Li et al., 2024), signaling a clear necessity for 103 methodological innovations in this area. 104

Addressing these concerns, we propose "Tree-of-Table", a novel paradigm crafted to optimize LLMs for the task of large-scale table understanding, as shown in Figure 1. By condensing and decomposing tables, our approach distills and systematizes the critical information into a tree-structured model that resonates with the stepwise reasoning employed by humans. This tree acts as a roadmap,

108 guiding the LLM through the complexities of the table in a logical and organized manner. It 109 provides a structured approach where each node serves a purpose, simplifying the interaction between 110 the LLM and the tabular data. The efficacy of our Tree-of-Table methodology is emphatically 111 validated through rigorous testing across a selection of datasets (including WikiTQ (Pasupat & Liang, 112 2015), TableFact (Chen et al., 2019), FeTaQA (Nan et al., 2022), and BIRD (Li et al., 2024)) with each presenting its own unique challenges. Consistently achieving top-tier results, Tree-of-Table 113 demonstrates not just its capacity to navigate the intricacies of table reasoning but also its potential to 114 set a new benchmark in the field. 115

116 117

118

2 RELATED WORK

Traditional Table Understanding. At the core, traditional methods have focused on generating
executable languages like SQL (Rajkumar et al., 2022); Liu et al., 2021; Eisenschlos et al., 2020;
Jiang et al., 2022) to interact with tables. This approach stems from the need to reason over both
free-form natural language questions and (semi-)structured tables. While effective in accessing tabular
data, these methods often fall short in capturing the nuanced semantics within a table, particularly
struggling with web tables that feature free-form text in cells.

Prompting Language Models for Table Understanding. A novel stride in table understanding has 126 been the application of prompting strategies (Wei et al., 2022; Chen et al., 2022; OpenAI, 2023; Imani 127 et al., 2023; Khot et al., 2022; Zhang et al., 2022). By generating reasoning steps through in-context 128 learning, models like Chain-of-Thought (Wei et al., 2022) and its evolutions (Yao et al., 2023a) break 129 down questions into sub-problems, iteratively solving each to improve comprehension of complex 130 tasks. These methods showcase LLMs' prowess in handling intricate reasoning chains, albeit not being 131 explicitly designed for tabular data. Emerging approaches have sought to extend LLM capabilities 132 beyond text, incorporating external tools to solve reasoning tasks (Cheng et al., 2022; Hsieh et al., 133 2023; Dhingra et al., 2019; Liu et al., 2023). Generating Python or SQL programs (Cheng et al., 2022; 134 Nahid & Rafiei, 2024; Shi et al., 2020; Pönighaus, 1995) and executing them with interpreters or APIs 135 has shown promise in enhancing arithmetic and table-based reasoning. However, the performance of these program-aided methods sometimes falters in complex table scenarios due to the static 136 nature of tables in the reasoning process. Dater (Ye et al., 2023) dynamically modifies the tabular 137 context to aid in solving table-based tasks, albeit primarily focusing on data pre-processing with 138 limited operations. Subsequently, the Chain-of-Table (Wang et al., 2024) method is inspired by the 139 chain-of-thought (Wei et al., 2022) principle and performs implicit sub-table extraction. Contrarily, 140 our proposed Tree-of-Table approach is inspired by the tree-of-thought principle (Yao et al., 2023a), 141 creating adaptive tree-based reasoning chains that exploit the planning capabilities of LLMs for more 142 nuanced and context-specific table reasoning. 143

Table Understanding Datasets. Datasets like WikiTQ (Pasupat & Liang, 2015) and TableFact (Chen et al., 2019) have been instrumental in developing table understanding methods. These standard benchmarks provide a foundation but are often limited in size and complexity. BIRD (Li et al., 2024) represents a significant leap forward in the field, being one of the largest and most intricate datasets designed for table understanding to date. Spanning across 37 professional domains with a substantial size of 33.4 GB, BIRD offers over 12,000 examples gleaned from real-world databases. Its development involved modifying open-source relational databases and curating additional ones, all complemented by crowdsourced natural language questions and corresponding SQL queries.

151 152

3 TREE-OF-TABLE: UNLEASHING THE POWER OF LLMS

153 154 155

156

3.1 FORMULATION OF LARGE-SCALE TABLE UNDERSTANDING

In the domain of table understanding, the core challenge lies in accurately interpreting and extracting information from tabular data in response to a given natural language query or statement. The essence of table understanding can be encapsulated as the task of mapping a natural language question or statement Q to a corresponding output S that accurately reflects the information contained within a table T. This table can be characterized by its structure, which includes rows and columns, with each cell representing a specific data point. Formally, a table T can be divided into headers H and



Figure 2: Illustration of the initial phases in the Tree-of-Table methodology, encompassing Table Condensation (the upper part), followed by Table-Tree Construction (the lower part). Starting with a large-scale input table, the process selectively condenses the data, emphasizing task-relevant information. Subsequently, the decomposed elements are methodically reorganized into a Table-Tree, a hierarchical structure designed to streamline and guide the subsequent reasoning process.

195 196 197

209 210

215

192

193

194

data values D, where each header in H corresponds to a column in the table and D represents the collective data points contained within these columns.

Furthermore, table understanding involves not just the direct interpretation of tables but also potentially requires external knowledge K and conversion of table data into a format that is amenable to computational models. This is especially relevant for tasks that involve complex reasoning or necessitate an understanding beyond the explicit table content, such as requiring background knowledge or contextual understanding to correctly interpret the question or the data.

In our experiments, we utilize a total of four datasets: WikiTQ (Pasupat & Liang, 2015), Tab-Fact (Chen et al., 2019), FeTaQA (Nan et al., 2022), and BIRD (Li et al., 2024). For WikiTQ, TabFact, and FeTaQA, there is no external knowledge; therefore, the table understanding problem can be defined as finding a function or model $f(\cdot, \theta)$ that satisfies

$$S = f(Q, \langle H, D \rangle | \theta), \tag{1}$$

where θ represents the model parameters. In contrast, for the BIRD dataset, there is external knowledge used to explain specific terms in the questions, allowing us to define the table understanding problem as

$$S = f(Q, \langle H, D \rangle, K|\theta), \tag{2}$$

where K denotes the external knowledge.

216 3.2 OVERVIEW

218 In this work, we introduce a novel approach named "Tree-of-Table" devised to address the challenge 219 of table reasoning within large-scale table understanding datasets (e.g., BIRD) and real-world applications, as shown in Figure 2 and Figure 3. Our methodology encompasses several steps 220 designed to simplify and enhance the reasoning capabilities of LLMs when confronted with large, 221 interconnected tables. First, we condense the tables based on the specific requirements of the query. 222 This process identifies relevant portions of the tables, thereby reducing the cognitive load on LLMs. We then apply a tree-based decomposition strategy to segment large tables into smaller, manageable 224 units, guided by the relationships among tables, such as foreign keys, and the structure of the query. 225 Next, we construct a "Table-Tree" by reorganizing the condensed information into a hierarchical 226 structure. Each node in this tree represents a logical block of information or a step in the reasoning 227 process, mirroring the cognitive approach of breaking down complex problems into simpler sub-228 problems. Finally, we perform a sequential traversal of the constructed Table-Tree to derive answers 229 to the queries. This systematic traversal ensures logical progression through each node, allowing 230 for the synthesis of information and insights gained from previous steps. The iterative nature of this reasoning process culminates in a well-informed conclusion. 231

- 232
- 233 234

242

250

3.3 TABLE CONDENSATION AND DECOMPOSITION

S

Addressing the significant challenges posed by large-scale relational tables to LLMs requires a nuanced understanding of the specific difficulties in question. These challenges primarily stem from two aspects: (1) The intricate foreign key relationships among multiple tables, which are commonly defined using SQL syntax, may not be readily interpretable by LLMs due to their complexity and the specialized knowledge required to understand relational database schemas. (2) The sheer size of the tables often exceeds the input context limit of LLMs, making it impossible for these models to process the entirety of the data directly. To effectively address these issues, our methodology incorporates two key processes: Table Condensation and Tree-based Decomposition.

243 3.3.1 TABLE CONDENSATION

As shown in the upper part of Figure 2, our initial step involves condensing the tables based on the context of the question Q and any additional evidence provided. Since there are possibly multiple tables, this process employs LLMs to identify one sub-table relevant to Q from them through schemalinking (Lei et al., 2020). Following the identification of relevant schemas/headers, we merge these multiple tables to reduce redundancy and decrease their size,

$$subTable_Q = f(Schema_Link(Q, \{H\}, K)|\theta),$$
(3)

251 where $\{H\}$ indicates the headers of the tables. This condensation aims to recall one sub-table 252 pertinent to Q and eliminate superfluous table information, thereby enhancing the information density 253 related to Q within the tables and making it more manageable for LLM processing. Through a detailed analysis of the BIRD dataset, we observed that over 70% of questions involved tables whose length 254 exceeded the input limitations of current LLMs. Furthermore, more than 90% of these questions 255 pertained to at least two tables, with 20% involving four or more tables, significantly complicating 256 the understanding process for LLMs. Post-condensation, we found that the length of tables involved 257 in more than 60% of long questions was reduced below the LLM input limit, and all questions were 258 associated with a singular, condensed table, as shown in the upper part in Figure 2. 259

260 3.3.2 TREE-BASED DECOMPOSITION

Even with reduced size and complexity post-condensation, the tables might still be too lengthy or
intricate for LLMs to handle efficiently, occasionally still surpassing the models' input limits. To
mitigate this, we begin by breaking down the question Q into its most general components, delineating
the entire problem-solving process into several independent yet sequentially connected steps.

$$S = \mathcal{P}_{\text{decomp}}(\{S_i^1\}, r^1 | Q), \quad i < \text{MAXDegree},$$
(4)

where S is the final solution, $\{S_i^1\}$ is the firstly decomposed intermediate sub-solutions towards S, r^1 is the possible relationship between $\{S_i^1\}$. \mathcal{P}_{decomp} is the "Thought Decomposition Prompt". "MAXDegree" is the pre-defined maximum degree of the Table-Tree. This decomposition involves mapping out the key stages $(\{S_i^1\} \text{ and } r^1)$ of reasoning required to address Q. By doing so, we transform a potentially overwhelming task into a series of manageable sub-tasks, each contributing incrementally to the formulation of the final answer. $\{S_i^1\}$ and r^1 also serve as the root node of the first-level subtree we will construct in the Table-Tree, for example, as shown the lower part in Figure 2.

275

276 3.4 TABLE-TREE CONSTRUCTION277

Drawing inspiration from the Tree-of-Thought concept (Yao et al., 2023a), our table-tree structure
closely resembles how humans naturally approach problem-solving. The illustration of overall
construction is showed in the lower part in Figure 2. When faced with a complex problem, people
typically employ a "breadth-first" strategy: deconstructing the problem into several general, independent yet interconnected subprocesses and then iteratively refining each subprocess into finer-grained
solutions.

3.4.1 BREADTH-FIRST THOUGHT GENERATION

Within this framework, we utilize in-context learning to instruct LLMs on dynamically generating thoughts for the question in a breadth-first way. Based on the firstly decomposed $\{S_i^1\}$ and r^1 in Eq. 4, the following breadth-first thought generation process can be formulated as,

 $S = \mathcal{P}_{\text{decomp}}(\{S_i^1\}, r^1 | Q), \quad i < \text{MAXDegree},$

.

288 289 290

284

285 286

287

294

$$S_{i}^{1} = \mathcal{P}_{\text{decomp}}(\{S_{i,j}^{2}\}, r_{j}^{2}), \quad j < \text{MAXDegree},$$

$$\dots$$

$$S_{i,j,\dots}^{d} = \mathcal{P}_{\text{decomp}}(\{S_{i,j,\dots,k}^{d+1}\}, r_{i,j,\dots,k}^{d+1}), \quad k < \text{MAXDegree},$$
(5)

295 296 297

298

$$S^{d_{\max}-1}_{i,j,\dots,k,\dots} = \mathcal{P}_{\operatorname{decomp}}(\{S^{d_{\max}}_{i,j,\dots,k,\dots,l}\}, r^{d_{\max}}_{i,j,\dots,k,\dots,l}), \quad d_{\max} <= \operatorname{MAXDepth},$$

where S is the final solution, $\{S_i^1\}$ is the firstly decomposed intermediate solutions towards S, r^1 is the possible relationship between $\{S_i^1\}$. $\{S_{i,j,\dots,k}^{d+1}\}$, $r_{i,j,\dots,k}^{d+1}$, and so on. Note that $r_{i,j,\dots,k}^{d+1}$ may be empty in the actual decomposition process if we need not to consider the relationship between $\{S_{i,j,\dots,k}^{d+1}\}$. *d* is the depth of thought. "MAXDegree" and "MAXDepth" are the pre-defined maximum degree and depth of the Table-Tree, respectively.

To prevent excessive decomposition of thoughts that could lead to redundant or erroneous reasoning processes, we set a maximum value for the depth *d*, denoted as "MAXDepth", and follow (Yao et al., 2023a) to utilize the LM to deliberately reason about end thought states $\{S_{i,j,\ldots,k,\ldots,l}^{d_{\max}}\}, r_{i,j,\ldots,k,\ldots,l}^{d_{\max}}\}$. Such a deliberate heuristic can be more flexible than programmed rules.

309 3.4.2 ITERATIVE CONSTRUCTION 310

Building upon the breadth-first thought generation, we construct the Table-Tree by iteratively constructing child nodes level by level for $\{S_i^1\}$ and r^1 , until we reach the leaf nodes at the bottom of the tree. Each sub-thought corresponding to $\{S_{i,j,\ldots,k}^{d+1}\}$ and $r_{i,j,\ldots,k}^{d+1}$ is regard as intermediate node, which act as "thinking node" representing a subprocess that can be further decomposed. The end thought state $\{S_{i,j,\ldots,k,\ldots,l}^{d_{\max}}\}$, $r_{i,j,\ldots,k,\ldots,l}^{d_{\max}}$ are set to leaf nodes, which functions as "execution nodes". In concrete, leaf nodes are the actionable endpoints of Table-Tree where specific operations such as data retrieval, calculations, or logical evaluations occur based on the parameters defined by their parent nodes. Therefore, we formulate $\{S_{i,j,\ldots,k,\ldots,l}^{d_{\max}}\}$, $r_{i,j,\ldots,k,\ldots,l}^{d_{\max}}$ as

319

$$S_{i,j,\dots,k,\dots,l}^{d_{\max}}, r_{i,j,\dots,k,\dots,l}^{d_{\max}} = \mathcal{P}_{\text{sample}}(\{\text{OP_Pool}\}), \quad d_{max} <= \text{MAXDepth.}$$
(6)

where \mathcal{P}_{sample} is the "Operation Sample Prompt". In selecting the operation pool, we based it on (Wang et al., 2024) and chose the most frequently used table operations from the resource at (Bytescout, 2024).



Figure 3: Depiction of the Table-Tree Execution phase within the Tree-of-Table approach. The model traverses the hierarchical Table-Tree, processing each node sequentially from the root to the leaves. At each step, the model integrates the information from the current node with the insights gathered from previous nodes, systematically building upon the reasoning chain to derive the final answer.

348 349 350

351

345

346

347

3.5 TABLE-TREE EXECUTION

After constructing the Table-Tree, we view it as a proxy task for the entire table understanding procedure. As shown in Figure 3, by traversing and executing operations across this tree, LLMs can implicitly generate tables and save intermediary results, thus enabling a seamless reasoning process. This stage diverges from the construction phase by utilizing a depth-first search approach to execute the thought chain, ensuring a systematic and comprehensive exploration of the tree structure.

357
 358
 359
 360
 360
 361
 362
 362
 364
 365
 365
 366
 366
 367
 368
 369
 360
 360
 361
 362
 362
 363
 364
 365
 365
 366
 366
 367
 368
 368
 369
 360
 360
 361
 361
 362
 362
 362
 363
 364
 364
 365
 365
 365
 366
 366
 366
 367
 368
 368
 369
 369
 360
 360
 361
 361
 362
 362
 362
 362
 363
 364
 364
 365
 365
 366
 366
 366
 367
 368
 368
 368
 369
 369
 360
 360
 361
 362
 362
 362
 362
 362
 362
 362
 363
 364
 364
 365
 365
 366
 366
 367
 368
 368
 368
 369
 369
 369
 360
 360
 361
 362
 362
 362
 362
 362
 364
 365
 365
 366
 366
 366
 367
 368
 368
 368
 368
 368
 368
 369

Leveraging Tree Structure for Efficiency. A key advantage of the tree structure is its inherent ability to enhance reasoning efficiency by logically organizing and compartmentalizing different aspects of the problem-solving process into subtrees. To exploit this benefit to its fullest, we execute the reasoning process subtree by subtree, based on the root node's children. After processing a subtree, we store its result before proceeding. This approach contrasts with linearly merging all subtrees into a single chain for execution. By maintaining the distinction between subtrees and executing them as separate units, we significantly mitigate the risk of intermediary tables becoming excessively large and unwieldy, which in turn, would thwart the reasoning process.

371

373

372 3.6 COMPARATION WITH CHAIN-OF-TABLE

In summary, the differences between the two methods are as follows: (1) Planning Style: Tree-ofTable uses a hierarchical, tree-based thought decomposition approach, while Chain-of-Table uses a
linear thought decomposition approach. (2) Execution Strategy: Notably, Chain-of-Table processes
the entire chain of history for each dynamic planning step and is shown to be effective for relatively
small tables, such as WikiTQ. However, as tables grow larger and questions become more complex,

~	-		
	. /	5	ζ.
J		~	ρ.

381 382

391 392 393

397

399

400

401

402

403

404

405

406

Table 1: Comparison of Table Understanding results on WikiTQ, TabFact datasets, with GPT3.5, PaLM2 and LLaMA2.

Mathad		WikiTQ		TabFact			
Method	GPT3.5	PaLM2	LLaMA2	GPT3.5	PaLM2	LLaMA2	
Text-to-SQL (Rajkumar et al., 2022b)	52.90	52.42	36.14	64.71	68.37	64.03	
End-to-End QA (Wang et al., 2024)	51.84	60.59	23.90	70.45	77.92	44.86	
Few-Shot QA (Wang et al., 2024)	52.56	60.33	35.52	71.54	78.06	62.01	
Binder (Cheng et al., 2022)	56.74	54.88	30.92	79.17	76.98	62.76	
Chain-of-Thought (Wang et al., 2024)	53.48	60.43	36.05	65.37	79.05	60.52	
Dater (Ye et al., 2023)	52.81	61.48	41.44	78.01	84.63	65.12	
Chain-of-Table (Wang et al., 2024)	59.94	67.31	42.61	80.20	86.61	67.24	
TREE-OF-TABLE	61.11	68.77	44.01	81.92	87.88	69.33	

Table 2: Comparison of Table Understanding results on FetaQA and BIRD datasets.

	FeTaQA			BIRD				
Method	BLEU	ROUGE-1	ROUGE-2	ROUGE-L	BLEU	ROUGE-1	ROUGE-2	ROUGE-L
FT(T5-large) (Ye et al., 2023)	30.54	0.63	0.41	0.53	-	-	-	-
End-to-End QA (Wang et al., 2024)	28.37	0.63	0.41	0.53	9.90	0.44	0.18	0.43
Codex (Chen et al., 2021)	27.96	0.62	0.40	0.52	-	-	-	-
Dater (Ye et al., 2023)	29.47	0.63	0.41	0.53	10.65	0.44	0.18	0.43
Chain-of-Table (Wang et al., 2024)	32.61	0.66	0.44	0.56	12.12	0.49	0.22	0.48
TREE-OF-TABLE	34.73	0.68	0.46	0.58	15.70	0.53	0.26	0.52

maintaining the complete thought chain becomes cumbersome, ultimately decreasing the efficiency of the model. Our Tree-of-Table method addresses this by embracing the tree's inherent "divide and conquer" philosophy (Bentley, 1980) to construct a Table-Tree. Each generation of child nodes relies exclusively on the information from their multi-level parent nodes, without the need for uncle nodes, as illustrated in Figure 2. By this way, we significantly reduce the historical chain's length on which each node's dynamic planning relies, to less than the depth of the tree, thus considerably simplifying the generation process at each level. (3) Data Preprocessing: Chain-of-Table does not preprocess complex tables with rich foreign key connections, whereas Tree-of-Table includes a Table Condensation step to handle such complexities.

407 408 409

410 411

412

4 EXPERIMENTS

4.1 DATASETS, METRICS, IMPLEMENTATION DETAILS

413 We evaluate our method on both three small table understanding benchmarks: WikiTQ (Pasupat & 414 Liang, 2015), FeTaQA (Nan et al., 2022), and TabFact (Chen et al., 2019), and one large-scale dataset: 415 BIRD (Li et al., 2024). For WikiTQ and TabFact, we employ the standard denotation accuracy metric. The nature of FeTaQA and BIRD for requiring elaborate responses prompts us to assess performance 416 through a variety of metrics including BLEU, ROUGE-1, ROUGE-2, and ROUGE-L (Lin, 2004) to 417 capture different facets of response quality. For our experiments, following the previous works (Wang 418 et al., 2024), we leverage the computational prowess of advanced language models, namely PaLM 419 2 (Anil et al., 2023), GPT 3.5 (OpenAI, 2023), and LLaMA 2 (Touvron et al., 2023). To facilitate 420 in-context learning, we incorporate a few-shot approach using demo samples from the training set 421 within the prompts, ensuring the models can effectively learn from limited examples.

422 423

4.2 MAIN RESULTS

In our experimental evaluation, we comprehensively compare our proposed approach, Tree-of-Table,
with several renowned baselines and state-of-the-art methodologies across both small and largescale tabular datasets, including WikiTQ, TableFact, FeTaQA, and BIRD, in Table 1 and Table
Our analysis is designed to assess the effectiveness of Tree-of-Table in facilitating complex
table understanding and reasoning tasks, particularly highlighting its performance in challenging
scenarios involving large-scale tables. The results in Table 1, demonstrate that Tree-of-Table not only
significantly outperforms early Generic Reasoning methods (End-to-End QA, Few-Shot QA, Chainof-Thought) and Program-aid Reasoning methods (Text-to-SQL, Dater, Binder) but also surpasses



Figure 4: Ablation study: (a) Generalization Ability under Different Table Sizes. (b) Effectiveness of Table Condensation.

the current state-of-the-art method, Chain-of-Table. This superiority was consistent across multiple LLMs including GPT3.5, PalM2, and LLAMA2. Our Tree-of-Table approach showcased robust enhancements in reasoning over both small and large tables. Specifically, in datasets with larger tables like BIRD, the benefits of Tree-of-Table were even more pronounced, suggesting that our tree-based method is particularly suited for complex, large-scale reasoning tasks where effective condensation, decomposition, and subtree execution strategies are critical.

4.3 ABLATION STUDY

444

445

446 447

448

449

450

451

452

453 454

455

456 Generalization Ability under Different Table Sizes. Here, we evaluate the generalization ca-457 pacity of our method across tables of varying sizes. Large tables present considerable compre-458 hension challenges to models, as their capacity to grapple with extended contexts and prompts 459 is severely tested. In Figure 4, we provide a detailed comparison of 4 methodologies (Binder, 460 Dater, Chain-of-Table, and Tree-of-Table) across two datasets (WikiTQ and BIRD). The per-461 formance metrics in table understanding tasks evidently deteriorate as the size of the tables increases. This degradation in performance reflects the inherent difficulties associated with 462 large table comprehension, confirming that it remains an exceptionally challenging problem area. 463 However, it is notable that the decline

in performance with increasing table

- 465 size is much more gradual for Tree-of-
- Table as compared to other methods.
- 467 Especially on the large-scale table
- dataset BIRD, Tree-of-Table demon strates superior robustness and gener-
- 470 alization ability. Even as table size
- scaled up, Tree-of-Table maintaines alevel of performance that was not only

Table 3: The	Node Number	and Height	of Tree Chains.
14010 01 1110	1.0000	and rivight	or rive chambr

Dataset	WikiTQ	TabFact	BIRD
Chain-of-Table: Chain Length	4	4	11
Tree-of-Table: Tree Height	3	3	7
Tree-of-Table: Node Number	6	7	18

better than its counterparts but also displayed less variance in its results.

474 Node Number and Height of Tree Chains. The height of the Table-Tree reflects the depth of the 475 model's reasoning chain, indicative of the complexity of the reasoning process. In contrast, the node 476 number corresponds to the length of the thought chain, representing the number of discrete reasoning 477 steps taken by the model. In Table 3, our comparative analysis reveals that across both smaller and larger datasets, the average height of the Table-Trees is generally less than that of the Chain-of-Table. 478 This indicates that the reasoning process in the Tree-of-Table method tends to require fewer levels 479 of hierarchical reasoning to arrive at a solution. Additionally, the average length of the Table-Trees 480 remains within a reasonable range, suggesting a good balance between depth and breadth in our tree 481 structures. 482

Comparison of Table Format Encoding. The encoding format of tables plays a vital role in how effectively a model can interpret and manipulate table data. Early research has indicated that the specific form of table encoding can significantly impact the model's performance in table understanding tasks. Here, we follow the lead of prior work, comparing the effects of four distinct

486 encoding formats on the final performance of table understanding: PIPE, HTML, TSV (Tab Separated 487 Values), and Markdown. As shown in Table 4, the Markdown format leads to the highest performance 488 among the tested encoding styles. The benefits of Markdown may be likely attributed to its readability, 489 clear structure, and straightforward syntax, all of which align well with the parsing capabilities of 490 LLMs.

491 Efficiency Analysis. In this context, efficiency refers to the model's 492 capability to achieve its goals with the least amount of computational 493 resources—specifically, the number of samples it needs to generate 494 to arrive at a correct answer. To substantiate the efficiency of our 495 Tree-of-Table methodology, we scrutinize how it compares with ex-496 isting methods in terms of the number of required generated samples to solve tasks. For a comprehensive analysis, we compared Tree-497 of-Table against notable methods such as Dater and Chain-of-Table 498 on BIRD. As depicted in Table 5, our analysis demonstrates that 499 Tree-of-Table consistently requires the fewest generated samples to 500 reach accurate answers across all evaluated datasets. 501

Table 4: Comparison of Table Format Encoding

onnat Encounty.	
Table Formatting	WikiTQ
HTML	68.01
TSV	68.12
PIPE	69.34
MarkDown	69.77

Method	Generate Samples
Dater	300
Chain-of-Table	120
TREE-OF-TBALE	90

Effectiveness of Table Condensation. Finaly, we validate the efficacy of the proposed Table Condensation component in reducing table sizes, making them more amenable to LLMs for reasoning. By condensing tables, we aim to filter out irrelevant information, thereby boosting the signal-to-noise ratio and allowing the model to focus on the most pertinent data. We conducte a comparative analysis of the number of table cells before and after applying Table Condensation across four datasets: WikiTQ, Tab-Fact, FeTaQA, and BIRD. Figure 4 (b) highlights the stark

511 contrast in table sizes before and after the application of Table Condensation. Across all examined 512 datasets, there is a significant reduction in the number of table cells post-condensation. This reduction 513 demonstrates the effectiveness of our method in shrinking table dimensions, ensuring that tables 514 remain within a tractable size range and contain information that is highly relevant to the task at hand. 515

5 CONCLUSION

517 518

521

516

504

505

506

507

508

509

510

519 In this paper, we address the profound challenge of advancing table understanding with LLMs, 520 specifically in the domain of large and complex tabular datasets. Our innovative approach, Treeof-Table, integrates table condensation and decomposition with a hierarchical reasoning construct 522 that aligns with human cognitive processes to tackle intricate problem-solving tasks. Our extensive 523 experiments conduct across various datasets, including WikiTQ, TableFact, FeTaQA, and BIRD, 524 demonstrate that Tree-of-Table not only achieves state-of-the-art performance but also presents 525 remarkable improvements in efficiency and generalizability.

526 Limitations. While Tree-of-Table has shown exceptional results in enhancing the efficiency and 527 effectiveness of large language models in processing extensive tabular data, the tree-based reason-528 ing may need to require careful calibration to balance depth and breadth effectively-a task that 529 necessitates fine-tuning and may impose certain limitations on adaptability.

530 **Broader Impact.** The broader impact of Tree-of-Table is multi-faceted, extending across academic, 531 industrial, and societal domains. Academically, our work contributes a significant leap forward in 532 the intersection of table understanding and natural language processing, providing a reference point 533 for future research and development in this area. In industry, the application of Tree-of-Table can 534 revolutionize the way organizations interact with large datasets. By simplifying the complexity and enhancing the reasoning capabilities of models with tabular data, Tree-of-Table can facilitate more 536 informed decision-making, enhance prediction systems, and optimize data-driven strategies across 537 various sectors such as finance, healthcare, and logistics. From a societal perspective, improving the accessibility and comprehension of large-scale data has the potential to democratize information. 538 By enabling a more nuanced understanding of data presented in tabular form, Tree-of-Table can contribute to greater transparency and empower individuals to make better data-informed decisions.

540 REFERENCES 541

542 543 544	Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, et al. Chain-of-table: Evolving tables in the reasoning chain for table understanding. <i>arXiv preprint arXiv:2401.04398</i> , 2024.
545 546 547	Michael J. Cafarella, Alon Halevy, Daisy Zhe Wang, Eugene Wu, and Yang Zhang. Webtables: Exploring the power of tables on the web. <i>Proc. VLDB Endow.</i> , 1(1):538–549, aug 2008. ISSN 2150-8097. doi: 10.14778/1453856.1453916.
549	OpenAI. Gpt-4 technical report. ArXiv, abs/2303.08774, 2023.
550 551 552	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of Thoughts: Deliberate problem solving with large language models, 2023b.
553 554	Wenhu Chen. Large language models are few (1)-shot table reasoners. In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pp. 1120–1130, 2023.
555 556 557 558	Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. Structgpt: A general framework for large language model to reason over structured data. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 9237–9251, 2023.
559 560 561 562	Shima Imani, Liang Du, and Harsh Shrivastava. MathPrompter: Mathematical reasoning using large language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track), pp. 37–42, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-industry.4.
563 564 565 566	Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. <i>arXiv</i> preprint arXiv:2305.10403, 2023.
567 568 569	Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Large language models still can't plan (a benchmark for llms on planning and reasoning about change). In <i>NeurIPS 2022 Foundation Models for Decision Making Workshop</i> , 2022.
570 571	Md Mahadi Hasan Nahid and Davood Rafiei. Tabsqlify: Enhancing reasoning capabilities of llms through table decomposition. <i>arXiv preprint arXiv:2404.10150</i> , 2024.
572 573 574 575 576	Si-An Chen, Lesly Miculicich, Julian Martin Eisenschlos, Zifeng Wang, Zilong Wang, Yanfei Chen, Yasuhisa Fujii, Hsuan-Tien Lin, Chen-Yu Lee, and Tomas Pfister. TableRAG: Million-token table understanding with language models. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> , 2024.
577 578 579 580	Yuan Sui, Jiaru Zou, Mengyu Zhou, Xinyi He, Lun Du, Shi Han, and Dongmei Zhang. TAP4LLM: Table provider on sampling, augmenting, and packing semi-structured data for large language model reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , 2024b.
581 582 583 584 585	Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets llm: Can large language models understand structured table data? a benchmark and empirical study. In <i>Proceedings of the 17th ACM International Conference on Web Search and Data Mining</i> , pp. 645–654, 2024a.
586 587 588	Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. Large language models are versatile decomposers: Decompose evidence and questions for table-based reasoning. <i>arXiv</i> preprint arXiv:2301.13808, 2023.
589 590 591	Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, et al. Binding language models in symbolic languages. In <i>International Conference on Learning Representations</i> , 2022.
592 593	Ziqi Jin and Wei Lu. Tab-cot: Zero-shot tabular chain of thought. <i>arXiv preprint arXiv:2305.17812</i> , 2023.

594 595 596	Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. In <i>International Conference on Learning Representations</i> , 2019.
597 598 599 600	Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
601 602	Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. Evaluating the text-to-sql capabilities of large language models. <i>arXiv preprint arXiv:2204.00498</i> , 2022b.
603 604 605	Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. Evaluating the text-to-sql capabilities of large language models. <i>arXiv preprint arXiv:2204.00498</i> , 2022a.
606 607 608	Tianze Shi, Chen Zhao, Jordan Boyd-Graber, Hal Daumé III, and Lillian Lee. On the potential of lexico-logical alignments for semantic parsing to sql queries. <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , 2020.
609 610 611 612	Richard Pönighaus. 'favourite'sql-statements—an empirical analysis of sql-usage in commercial applications. In <i>International Conference on Information Systems and Management of Data</i> , pp. 75–91. Springer, 1995.
613 614	George Katsogiannis-Meimarakis and Georgia Koutrika. A survey on deep learning approaches for text-to-sql. <i>The VLDB Journal</i> , pp. 1–32, 2023.
615 616 617 618	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837, 2022.
619 620 621 622 623	Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics</i> <i>and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long</i> <i>Papers)</i> , pp. 1470–1480, Beijing, China, July 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1142.
624 625 626 627 628	Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Hailey Schoelkopf, Riley Kong, Xiangru Tang, Mutethia Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, Dragomir Radev, and Dragomir Radev. FeTaQA: Free-form table question answering. <i>Transactions of the Association for Computational Linguistics</i> , 10:35–49, 2022. doi: 10.1162/tacl_a_00446.
629 630 631 632	Qian Liu, Bei Chen, Jiaqi Guo, Morteza Ziyadi, Zeqi Lin, Weizhu Chen, and Jian-Guang Lou. TAPEX: Table pre-training via learning a neural sql executor. In <i>International Conference on Learning Representations</i> , 2021.
633 634 635 636	Julian Eisenschlos, Syrine Krichene, and Thomas Müller. Understanding tables with intermediate pre-training. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pp. 281–296, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.27.
637 638 639 640 641	Zhengbao Jiang, Yi Mao, Pengcheng He, Graham Neubig, and Weizhu Chen. OmniTab: Pretraining with natural and synthetic data for few-shot table-based question answering. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pp. 932–942, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.68.
642 643 644 645	Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompt- ing: Disentangling computation from reasoning for numerical reasoning tasks. <i>arXiv preprint</i> <i>arXiv:2211.12588</i> , 2022.
646 647	Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In <i>International Conference on Learning Representations</i> , 2022.

648 649 650	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In <i>International Conference on Learning Representations</i> , 2022.
651	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik
652	Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. arXiv
653	preprint arXiv:2505.10001, 2025a.
654	Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner,
655	Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger
656 657	language models with less training data and smaller model sizes. In <i>Findings of the Association</i> for <i>Computational Linguistics: ACL 2023</i> Association for Computational Linguistics 2023
658	for computational Emgassies. The 2020 The sociation for computational Emgassies, 2020.
659	Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William
660	Cohen. Handling divergent reference texts when evaluating table-to-text generation. In <i>Proceedings</i>
661	of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 4884–4895, 2019
662	2019.
663	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,
664	and Percy Liang. Lost in the middle: How language models use long contexts. arXiv preprint
665	arXiv:2307.03172, 2023.
666	Wengiang Lei, Weixin Wang, Zhixin Ma, Tian Gan, Wei Lu, Min-Yen Kan, and Tat-Seng Chua.
667	Re-examining the role of schema linking in text-to-sql. In <i>Proceedings of the 2020 Conference on</i>
668	Empirical Methods in Natural Language Processing (EMNLP), pp. 6943–6954, 2020.
669	Pritagenit https://hytogenit.com/hlag/20 important cal quaries html. In Putagenit 2024
670	Bytescout. https://bytescout.com/biog/20-important-sqi-queries.html. in <i>Bytescout</i> , 2024.
671	Robert Tarjan. Depth-first search and linear graph algorithms. SIAM journal on computing, 1(2):
672	146–160, 1972.
673	Ion Louis Bentley Multidimensional divide-and-conquer Communications of the ACM 23(A):
674 675	214–229, 1980.
676	Chin-Yew Lin ROUGE: A package for automatic evaluation of summarizes. In Text Summarization
677	Branches Out, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
678	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amiad Almahairi, Vasmine Bahaei, Nikolay
679	Bashlykov, Soumva Batra, Prajiwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
681	and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
682	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
683	Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
684	language models trained on code. arXiv preprint arXiv:2107.03374, 2021.
685	Mohammadreza Pourreza and Davood Rafiei Din-sal: Decomposed in-context learning of text-to-sal
686	with self-correction. Advances in Neural Information Processing Systems, 36, 2024a.
687	
688	Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, Qian-Wen Zhang, Zhao Yan, and
600	Znoujun Li. Mac-sql: Multi-agent collaboration for text-to-sql. arXiv preprint arXiv:2312.11242, 2023
601	2023.
692	Mohammadreza Pourreza and Davood Rafiei. Dts-sql: Decomposed text-to-sql with small large
693	language models. arXiv preprint arXiv:2402.01117, 2024b.
694	Dongiun Lee, Choongwon Park, Jaehvuk Kim, and Heesoo Park Mcs-sol. Leveraging multiple
695	prompts and multiple-choice selection for text-to-sql generation. arXiv preprint arXiv:2405.07467.
696	2024.
697	Danny Zhou Nathanaal Sahärli La Hou Jacon Wai Nathan Sacha Verschi Wars Dala Sah
698	Claire Cui Olivier Bousquet Ouoc Le et al Least-to-most prompting enables complex reasoning
699	in large language models. arXiv preprint arXiv:2205.10625, 2022.
700	
701	Jon Louis Bentley. Multidimensional binary search trees used for associative searching. <i>Communica-</i> <i>tions of the ACM</i> , 18(9):509–517, 1975.

Robert D Blumofe and Charles E Leiserson. Scheduling multithreaded computations by work stealing. Journal of the ACM (JACM), 46(5):720-748, 1999. Wei-Yin Loh. Classification and regression trees. Wiley interdisciplinary reviews: data mining and knowledge discovery, 1(1):14–23, 2011. Allen Newell. Human problem solving. Upper Saddle River/Prentive Hall, 1972. Allen Newell, John C Shaw, and Herbert A Simon. Report on a general problem solving program. In IFIP congress, volume 256, pp. 64. Pittsburgh, PA, 1959. George A Miller. The magical number seven, plus or minus two: Some limits on our capacity for processing information. Psychological review, 63(2):81, 1956.

756 A APPENDIX

758 A.1 MORE EXPERIMENTS

760 A.1.1 ABLATION STUDY ON TIME COSTS

In addition to the existing efficiency analysis based on the number of samples (in Sec 3.3, Table 5), we further include a time cost analysis in Table 6, which shows the comparison of overall time costs between Tree-of-Table and the current state-of-the-art, Chain-of-Table.

Table 6: Comparison of overall time costs between Tree-of-Table and Chain-of-Table.

Method	Chain-of-Table	Tree-of-Table
Time Cost	5.7s	7.8s

Table 7: Accuracy evaluation results for Tree-of-Table in comparison with other Text-to-SQL methods.

774	Methods	DIN-SQL	MAC-SQI	DTS-	SQL	MCS-SQI	L Tree-of-Table
775	Accuracy (%)	50.72	57.56	55.	80	63.36	65.07
776		1					
777							
778		Tab	la 8. Salact	ion of M	A X Dent	h	
779						10	
780		MAXI	Jepth	6	8	10	
781		Accura	$acy(\%) \mid 5$	9.96 6	1.11 6	0.47	
782							
783							
784		Tab	le 9 [.] Select	ion of M	axDegre	e	
785				2	4	<u> </u>	
786		MAAI	Jegree	3	4	3	
787		Accura	$acy(\%) \mid 6$	0.03 6	1.11 6	0.91	
788							
789							
790		r	Table 10: E	rror Prop	agation		
791	Dom	ove Dercent(0	ເ) 5	Q	10	12	15
792	Kell		5	0	10	12	15
793	Acc	uracy(%)	61.09	61.00	60.87	60.22	59.45
794							

A.1.2 ACCURACY EVALUATION RESULTS

In addition to Table 2, we also show the accuracy evaluation results for Tree-of-Table in comparison with other Text-to-SQL methods such as DIN-SQL (Pourreza & Rafiei, 2024a), MAC-SQL (Wang et al., 2023), DTS-SQL (Pourreza & Rafiei, 2024b), and MCS-SQL (Lee et al., 2024), in Table 7.

801 A.1.3 SELECTION OF MAXDEGREE AND MAXDEPTH

In our experiments, we initially conduct preliminary experiments to roughly determine their value ranges. We then perform detailed hyperparameter experiments to identify the optimal values. The ablation study results on WikiTQ are listed in the Table 9 and Table 8. It can be seen that our proposed method is relatively robust with respect to these two hyperparameters.

- A.1.4 ERROR PROPAGATION
- 809 We also conduct an error propagation study on WikiTQ in Table 10. Specifically, we manually remove important information randomly during the table condensation step (at the root) and evaluate

810 the performance. The experimental results show that even when a certain amount of important 811 information is removed, the performance of Tree-of-Table does not significantly degrade and still 812 achieves relatively correct results, indicating the robustness of our method. 813

814 A.1.5 CASE STUDY

815

817

819

821

822 823

824

829

830

831

832

833

834

835

836

837

838

839

840

841

842

To more intuitively verify the superiority of our method, we have also included specific case compar-816 isons in Table 11, which encompass two typical complex logic problems: precise calculation and statistics. 818

- A.2 PROMPT TEMPLATE 820
 - The main structure of prompt template of Tree-of-Table is shown in Figure 5.
 - THEORETICAL ADVANTAGES OF HIERARCHICAL TREE-BASED ARCHITECTURE A.3

825 Previous work has shown that tree-based structures exhibit significant advantages in LLM reasoning, 826 as demonstrated in (Yao et al., 2023b; Zhou et al., 2022). In addition, earlier works also demonstrate 827 from multiple perspectives that hierarchical tree-based architecture is superior to linear architecture, including but not limited to the following points: 828

- Divide and Conquer Strategy (Bentley, 1975): Tree-based architectures leverage the divideand-conquer approach, which is a well-established algorithmic paradigm. This strategy breaks a problem into smaller subproblems, solves each subproblem independently, and combines their results. This can lead to more efficient processing, especially in complex, large-scale problems.
- Scalability (Blumofe & Leiserson, 1999): Trees can handle larger and more complex datasets more efficiently than linear structures. As data grows, the depth of a tree increases logarithmically rather than linearly, allowing for more scalable processing.
 - Improved Decision Making (Loh, 2011): In decision-making processes, tree structures can better model decision paths and outcomes, providing clearer insights into the reasoning behind decisions.
- Cognitive Alignment (Newell, 1972; Newell et al., 1959; Miller, 1956): Human reasoning often aligns more closely with hierarchical structures, which can make tree-based models more intuitive and easier to interpret.
- 843 844 845
- A.4 REPHRASE THE PROCESS OF CHAIN-OF-TABLE

846 Overall, Chain-of-Table processes table understanding based on the Chain-of-Thought LLM reasoning 847 approach, which employs an intuitive linear thought chain to decompose the problem. At each step 848 of decomposition, it selects operations from a predefined operations pool and generates intermediate 849 results for table processing. However, when dealing with complex, multi-branch logic, this linear 850 approach can produce lengthy and disorganized thought processes, making table reasoning chaotic 851 and prone to errors. Therefore, this linear approach is generally suited for relatively simple problems 852 and smaller tables.

853 In contrast, our proposed Tree-of-Table method is divided into three main parts: (1) Table Condensa-854 tion (Sec. 3.3): Given the input table, which contains many large-scale relevant tables connected via 855 foreign keys, we first employ LLMs to condense the table. This helps recall the sub-tables pertinent 856 to the query and eliminate extraneous information. (2) Table-Tree Construction (Sec. 3.4): Following the preprocessing steps, we instruct the LLM to dynamically generate child nodes in the tree. Each 858 generation of child nodes relies solely on their multi-level parent nodes, without needing reference to 859 uncle nodes. We then iteratively construct child nodes level by level until we reach the leaf nodes at the bottom of the tree. (3) Table-Tree Execution (Sec. 3.5): Finally, after constructing the Table-Tree, 860 we consider it a proxy task for the entire table understanding process. The traversal and execution of 861 operations across this tree enable a seamless reasoning process. 862

In summary, the differences between the two methods are as illustrated in Sec. 3.6.

Table 11:	Case	Study.
-----------	------	--------

880	Question	Chain-of-Table	Tree-of-Table	Groundtruth
881 882 883 884	"What was the growth rate of the total amount of loans across all accounts for a male client between 1996 and 1997?"	0	25.30	25.30
885 886 887 888 888 889 890 891	"Consider the average difference between K-12 enrollment and 15-17 enrollment of schools that are locally funded, list the names and DOC type of schools which has a difference above this average"	Incomplete list	Mountain Oaks(00), Castle Rock(00), Charter Community School Home Study Academy(00), Clovis Online Charter(54), Washington Elementary(52), 	Mountain Oaks(00), Castle Rock(00), Charter Community School Home Study Academy(00), Clovis Online Charter(54), Washington Elementary(52),
892 893 894 895 896 896 897 898 898	"What is the e-mail address of the administrator of the school located in the San Bernardino county, District of San Bernardino City Unified that opened between 1/1/2009 to 12/31/2010 whose school types are public Intermediate/Middle Schools and Unified Scools?"	www.realjourney.org	a.lucero@realjourney.org	a.lucero@realjourney.org
900 901 902 903 904	"For the branch which located in the south Bohemia with biggest number of inhabitants, what is the percentage of the male clients?"	40	44.26	44.26
905 906				

919		
920		
921		
922		
923	## Task Instruction ##	Yearmonth: [
924	You are an information and database analysis expert. When faced with a specific database with some tables and a question, please use logical analysis methods to	(CustomerID, customer ID. Value example: [3, 5, 6].), (Date, date, Value example: [2012-08-24].)
925	rigorously output the answer from the provided information, given the instruction.	(Consumption, consumption. Value example: [528.3])
926	## Guidelines ##] (foreign key: CustomerID)
927	Before answering the question, you will be presented with some /*Database tables*/,	/*Ougstion*/
928	You will rely on the provided information, utilizing methods of logical analysis, to	Which of the three segments—SME, LAM and KAM—has the
929	answer the question using only the provided materials. Aiming to provide an accurate answer.	biggest and lowest percentage increases in consumption paid in EUR between 2012 and 2013?
930	/*!*	/***:::*/
931	/*instruction*/ Perform Table-Tree Construction based on the operation pool and then Execute the	/*Evidence*/ Increase or Decrease = consumption for 2013 - consumption
932	Table-Tree.	for 2012; Percentage of Increase = (Increase or Decrease /
033	[Table-Tree Construction]:	be represented by Between 201201 And 201312; First 4
03/	- Thought Decomposition: You should generate a Tree-like thought chains for the	strings of Date represents the year.
025	question and tables. Following by the "divide and conquer" philosophy, each	/*Output of Table-Tree */
036	without the need for uncle nodes. The intermediate nodes represent a subprocess	(Tree-1 Chain): Increase Percentage of Consumption for SME
930	that can be further decomposed.	(i) \rightarrow Consumption of SME (i) \rightarrow EUR (i) \rightarrow 2012 (i) \rightarrow [sum sum cumulative] (l) \rightarrow 2013 (i) \rightarrow [sum sum cumulative]
90 <i>1</i> 000	- Operation Sample: The leaf nodes are the actionable endpoints of the tree where	(I) \rightarrow others (i) \rightarrow EUR (i) \rightarrow Not EUR (i) \rightarrow Consumption of
930	specific operations such as data retrieval, calculations can execute based on the parameters defined by their parent nodes. You will sample the operations from the	SME (i) \rightarrow Percentage of Increase of SME (i) \rightarrow Increase (i) \rightarrow subtract (I) \rightarrow Percentage (i) \rightarrow divide (I) \rightarrow Percentage of
939	Operation Pool	Increase of SME (i) \rightarrow Increase Percentage of Consumption
940	[Table-Tree Execution]: Execute the table-tree by depth-first searching	\rightarrow (Tree-2 Chain): // omission of Tree-2 Chain
941	/*Operation Pool*/	\rightarrow (Tree-3 Chain): // omission of Tree-3 Chain \rightarrow (Tree-4 Chain): Compare biggest and lowest (i) \rightarrow sort by
942	// omission of operation pool, following Chain-of-Table	(I) (I)
943	## Some In-context examples ##	The 4 chains are execute one by one
944	/*Database tables*/	/*Answer*/ KAM has the higgest nercentage increases SME has the
945	The condensated table is:	lowest percentage increase.
946	// omission of table content	# provide more examples
947	/*Database schema*/	# amission of anaration neal, could refer to Chain of
948	(CustomerID, customer ID. Value example: [3, 5, 6].),	Table
949	(Segment, consumption type. Value example: ['SME', 'LAM', 'KAM'].),	
950] (foreign key: CustomerID, Segment)	## Query ##
951	Gasstations: [/*Database tables*/ // omission of table content
952	(GasStationID, gas station ID. Value example: [44, 45, 60].),	/*Databasa schoma*/
953	(Country, country name. Value example: ['CZE']),	
954	(Segment, consumption type. Value example: ['SME', 'LAM', 'KAM'].)] (foreign kev: CustomerID. Segment, GasStationID)	(sname, school name. Value examples: [Making Waves Academv].).
955		(cname, county name. Value examples: [Contra Costa].),
956	(Products: [(ProductID, product ID. Value example: [1, 3, 5].),	examples: [73].),
957	(Description, product description. Value example: ['Special', 'Protraviny'].),]
958	j (oreign key. Housenb)	
959	Transactions_1k: [(TransactionID, transaction ID. Value example: [11, 13, 16].),	/*Question*/
960	(Date, date. Value example: [2012-08-24].),	Which school in Contra Costa has the highest number of test
961	(CustomerID, customer ID. Value example: [3, 5, 6].),	LakeiSf
962	(CardID, card ID. Value example: [486621].), (GasStationID, gas station ID. Value example: [44, 45, 60].).	/*Evidence*/ // omission
963	(ProductID, product ID. Value example: [1, 3, 5].),	
964	(Amount, amount. Value example: [28, 18].), (Price, price. Value example: [2, 8].),	/*Output of Table-Tree */
965] (foreign key: CustomerID, ProductID, GasStationID)	/*Answer*/
966		

Figure 5: The main structure of prompt template of Tree-of-Table.